

# Large-scale sequencing of SARS-CoV-2 genomes from one region allows detailed epidemiology and enables local outbreak management

Andrew J. Page<sup>1,\*</sup>, Alison E. Mather<sup>1,2</sup>, Thanh Le-Viet<sup>1</sup>, Emma J. Meader<sup>3</sup>, Nabil-Fareed Alikhan<sup>1</sup>, Gemma L. Kay<sup>1</sup>, Leonardo de Oliveira Martins<sup>1</sup>, Alp Aydin<sup>1</sup>, David J. Baker<sup>1</sup>, Alexander J. Trotter<sup>1,2</sup>, Steven Rudder<sup>1</sup>, Ana P. Tedim<sup>1,4</sup>, Anastasia Kolyva<sup>1,3</sup>, Rachael Stanley<sup>3</sup>, Muhammad Yasir<sup>1</sup>, Maria Diaz<sup>1</sup>, Will Potter<sup>3</sup>, Claire Stuart<sup>3</sup>, Lizzie Meadows<sup>1</sup>, Andrew Bell<sup>1</sup>, Ana Victoria Gutierrez<sup>1</sup>, Nicholas M. Thomson<sup>1</sup>, Evelien M. Adriaenssens<sup>1</sup>, Tracey Swingle<sup>2</sup>, Rachel A. J. Gilroy<sup>1</sup>, Luke Griffith<sup>2</sup>, Dheeraj K. Sethi<sup>3</sup>, Dinesh Aggarwal<sup>5,6,7,8</sup>, Colin S. Brown<sup>5</sup>, Rose K. Davidson<sup>2</sup>, Robert A. Kingsley<sup>1,2</sup>, Luke Bedford<sup>9</sup>, Lindsay J. Coupland<sup>3</sup>, Ian G. Charles<sup>1,2</sup>, Ngozi Elumogo<sup>1,3</sup>, John Wain<sup>1,2</sup>, Reenesh Prakash<sup>3</sup>, Mark A. Webber<sup>1,2</sup>, S. J. Louise Smith<sup>10</sup>, Meera Chand<sup>5</sup>, Samir Dervisevic<sup>3</sup>, Justin O'Grady<sup>1,2</sup> and The COVID-19 Genomics UK (COG-UK) Consortium

## Abstract

The COVID-19 pandemic has spread rapidly throughout the world. In the UK, the initial peak was in April 2020; in the county of Norfolk (UK) and surrounding areas, which has a stable, low-density population, over 3200 cases were reported between March and August 2020. As part of the activities of the national COVID-19 Genomics Consortium (COG-UK) we undertook whole genome sequencing of the SARS-CoV-2 genomes present in positive clinical samples from the Norfolk region. These samples were collected by four major hospitals, multiple minor hospitals, care facilities and community organizations within Norfolk and surrounding areas. We combined clinical metadata with the sequencing data from regional SARS-CoV-2 genomes to understand the origins, genetic variation, transmission and expansion (spread) of the virus within the region and provide context nationally. Data were fed back into the national effort for pandemic management, whilst simultaneously being used to assist local outbreak analyses. Overall, 1565 positive samples (172 per 100000 population) from 1376 cases were evaluated; for 140 cases between two and six samples were available providing longitudinal data. This represented 42.6% of all positive samples identified by hospital testing in the region and encompassed those with clinical need, and health and care workers and their families. In total, 1035 cases had genome sequences of sufficient quality to provide phylogenetic lineages. These genomes belonged to 26 distinct global lineages, indicating that there were multiple separate introductions into the region. Furthermore, 100 genetically distinct UK lineages were detected demonstrating local evolution, at a rate of ~2 SNPs per month, and multiple co-occurring lineages as the pandemic progressed. Our analysis: identified a discrete sublineage associated with six care facilities; found no evidence of reinfection in longitudinal samples; ruled out a nosocomial outbreak; identified 16 lineages in key workers which were not in patients, indicating infection control measures were effective; and found the D614G spike protein mutation which is linked to increased transmissibility dominates the samples and rapidly confirmed relatedness of cases in an outbreak at a food processing facility. The large-scale genome sequencing of SARS-CoV-2-positive samples has provided valuable additional data for public health epidemiology in the Norfolk region, and will continue to help identify and untangle hidden transmission chains as the pandemic evolves.

## DATA SUMMARY

Raw reads are deposited in the European Nucleotide Archive, and individual accession numbers are listed in Table S1 (available in the online version of this article). Consensus sequences are deposited with GISAID (Global Initiative on Sharing All Influenza Data), and individual accession numbers are listed in Table S1. The data can be visualized in MicroReact [https://](https://beta.microreact.org/project/8HsL3eyA16WsjtymhvbbEa-cog-uk-2021-03-18-uk-sars-cov-2)

[beta.microreact.org/project/8HsL3eyA16WsjtymhvbbEa-cog-uk-2021-03-18-uk-sars-cov-2](https://beta.microreact.org/project/8HsL3eyA16WsjtymhvbbEa-cog-uk-2021-03-18-uk-sars-cov-2). Bioinformatics pipelines used to process the data are available from <https://github.com/quadram-institute-bioscience/ncov2019-artic-nf/tree/qib> which were adapted from <https://github.com/connor-lab/ncov2019-artic-nf>. Additional metadata, trees and alignments are available from COG-UK <https://www.cogconsortium>.

uk. Supplementary Material for this article can be found on figshare: <https://figshare.com/s/41951bcd144df26c528a>.

## INTRODUCTION

In December 2019, a new coronavirus-related disease (COVID-19) was first reported in Wuhan, China [1]; the causal agent was identified as the novel human coronavirus, SARS-CoV-2. Since then, SARS-CoV-2 has spread globally leading to 120 million confirmed infections and 2.7 million deaths (as of 17 March 2021) [2]. Two risk factors are associated with higher mortality: sex, as males are at higher risk than females; and age, as older age groups are at substantially higher risk [3].

Whole genome sequencing provides high-resolution data that enable investigation of pathogen evolution and population structure [4]. When combined with robust epidemiological data, it is possible to gain insights into SARS-CoV-2 origins [5], transmission (both global [6] and local [7]) and responses to control measures [8]. Since the start of the pandemic, sequencing efforts and data sharing have facilitated tracking of the pandemic [9], identifying multiple independent virus introductions into different countries [6]. The ability to assign identifiers rapidly to groups of samples that are related is essential in public health, as demonstrated for influenza [10]. These identifiers can be formulated in different ways: from conserved sequences identified by multi-locus sequence typing [11], by assigning SNP addresses [12]; or, in the case of SARS-CoV-2, through the assignment of lineages [13].

The COVID-19 Genomics UK (COG-UK) consortium [14] is a UK-wide public health surveillance initiative comprising nearly 20 organizations from universities, research institutes and public health agencies that was created to generate and analyse large-scale SARS-CoV-2 sequencing datasets to understand virus evolution, transmission and spread in the UK. These data allow detailed insight into the course of the pandemic at the country, county and individual institution level. It was through large-scale analysis of SARS-CoV-2 genomes that evidence of a mutation (D614G) in the spike

## Impact Statement

Genomic lineages of SARS-CoV-2 can be used to track progression of the pandemic on an international scale. We undertook whole genome sequencing of the SARS-CoV-2 genomes present in positive clinical samples from one region. We combined clinical metadata sequencing data to understand the origins, genetic variation, transmission and expansion of the virus within the region and provide context nationally in the UK. In total, 42.6% of all positive samples identified by hospital testing were sequenced. The large-scale genome sequencing of SARS-CoV-2-positive samples has provided valuable additional data for public health epidemiology in the Norfolk region, and will continue to help identify and untangle hidden transmission chains as the pandemic evolves.

protein was revealed; it is likely that this mutation is responsible for increased transmissibility of the virus [15].

For the Norfolk region, we established a robust, rapid sequencing pipeline for SARS-CoV-2. Weekly sequencing data were fed back into the national effort for pandemic management, whilst simultaneously being used to assist local outbreak analyses. Here we describe the sequencing of genomes present in 1565 SARS-CoV-2 samples from 1376 cases, collected between March and August 2020. This represented 42.6% of all cases in the local population and included those with a clinical need, and key workers (such as healthcare, care and police) and their families. For context, at the end of the study period (27 August 2020), only five countries (UK, Australia, Spain, India and the USA) out of 103 countries had sequenced more SARS-CoV-2 genomes than had been sequenced in Norfolk for this paper. We used these data to investigate the genetic and epidemiological characteristics of the COVID-19 pandemic in the stable, low-density population of Norfolk and surrounding areas, UK. Our objectives were to use these sequence data to understand the evolution and spread of the virus locally, adding context

Received 16 December 2020; Accepted 19 April 2021; Published 29 June 2021

**Author affiliations:** <sup>1</sup>Quadram Institute Bioscience, Norwich Research Park, Norwich, NR4 7UQ, UK; <sup>2</sup>University of East Anglia, Norwich Research Park, Norwich, NR4 7TJ, UK; <sup>3</sup>Norfolk and Norwich University Hospital, Colney Lane, Norwich, NR4 7UY, UK; <sup>4</sup>Grupo de Investigación Biomédica en Sepsis - BioSepsis, Hospital Universitario Río Hortega/Instituto de Investigación Biomédica de Salamanca (IBSAL), Valladolid/Salamanca, Spain; <sup>5</sup>Public Health England, 61 Colindale Ave., London, NW9 5EQ, UK; <sup>6</sup>Department of Medicine, University of Cambridge, Cambridge, UK; <sup>7</sup>Cambridge University Hospital NHS Foundation Trust, Cambridge, UK; <sup>8</sup>Wellcome Sanger Institute, Hinxton, Cambridge, UK; <sup>9</sup>Ipswich Hospital, Heath Road, Ipswich, IP4 5PD, UK; <sup>10</sup>Public Health, County Hall, Martineau Lane, Norwich, NR1 2DH, UK.

\*Correspondence: Andrew J. Page, [andrew.page@quadram.ac.uk](mailto:andrew.page@quadram.ac.uk)

**Keywords:** SARS-CoV-2; sequencing; NGS; genome; genomic epidemiology; ARTIC.

**Abbreviations:** BCRC, Bob Champion Research and Education Building; COG-UK, COVID-19 Genomics Consortium; ENA, European Nucleotide Archive; GISAID, global initiative on sharing all influenza data; IPC, infection prevention and control; JPUH, The James Paget University Hospital; NNUH, Norfolk and Norwich University Hospital; PHE, Public Health England; QC, quality control; QEH, The Queen Elizabeth Hospital; QIB, Quadram Institute Bioscience; RLU, relative fluorescence units.

Raw sequence data are deposited in the European Nucleotide Archive and consensus sequences are deposited in GISAID. Accession numbers are listed in Table S1.

†These authors contributed equally to this work

**Data statement:** All supporting data, code and protocols have been provided within the article or through supplementary data files. Ten supplementary figures are available with the online version of this article.

to the national and global data, and to evaluate the role of rapid whole-genome sequencing for outbreak analysis in this setting.

Our analysis: identified a sublineage associated with six care facilities; found no evidence of reinfection in longitudinal samples; ruled out a nosocomial outbreak; identified 16 lineages in key workers which were not in patients, indicating infection control measures were effective; and found the D614G spike protein mutation which is linked to increased transmissibility dominates the samples and rapidly confirmed relatedness of cases in an outbreak at a food processing facility. This demonstrates the valuable role of large-scale genome sequencing of SARS-CoV-2 to inform surveillance and regional outbreak management.

## METHODS

### Samples

The clinical samples we used were initially collected passively for diagnostic testing with ethical approval from Public Health England (R and D ref. NR0195) and with sampling directed by government public health policy and local clinical need. Samples were taken at four large hospitals: Norfolk and Norwich University Hospital (NNUH) (1200 beds) in Norwich, Norfolk; The Queen Elizabeth Hospital (QEH) (500 beds) in Kings Lynn, Norfolk; The James Paget University Hospital (JPUH) (500 beds) in Great Yarmouth, Norfolk; and the Ipswich Hospital (550 beds) in Ipswich, Suffolk. Additional clinical samples that were included were collected at five smaller hospitals; by three community care organizations (representing dozens of care facilities and GP practices); and at drive-through testing facilities for healthcare workers, essential workers (such as police) and their families who live or work in Norfolk and the surrounding areas (Fig. S1).

### Samples and RNA extraction

Samples from cases with suspected SARS-CoV-2 were processed using five different diagnostic platforms over three laboratories on the Norwich Research Park: the Cytology Department and Microbiology Department, NNUH, Norwich, UK and the Bob Champion Research and Education Building (BCRE), University of East Anglia, Norwich, UK. Samples were primarily nasal/oropharyngeal swabs, although nasopharyngeal aspirates, bronchoalveolar lavage and sputum samples were also collected. The Cytology Department processed samples using the Roche Cobas 8800 SARS-CoV-2 system (<https://tinyurl.com/yy58t8sp>). The Microbiology Department processed samples using either the Hologic Panther Fusion System SARS-CoV-2 assay (<https://tinyurl.com/yye3m25p>) according to the manufacturer's instructions, the AusDiagnostics SARS-CoV-2, Influenza and RSV 8-well panel (<https://tinyurl.com/yyeh5y2w>) or Altona Diagnostics RealStar SARS-CoV-s RT-PCR Kit 1.0 (<https://altona-diagnostics.com/en/products/reagents-140/reagents/realstar-real-time-pcr-reagents/realstar-sars-cov-2-rt-pcr-kit-ruo.html>). RNA was extracted from swab samples in the

Microbiology Department using either the QIA Symphony (Qiagen) or AusDiagnostics MT-Prep (AusDiagnostics) instruments according to the manufacturer's instructions before being processed through the AusDiagnostics assay. In the BCRE, RNA was extracted using the MagMAX Viral/Pathogen II Nucleic Acid Isolation kit (Applied Biosystems) according to the manufacturer's instructions and the KingFisher Flex system (ThermoFisher). The presence of SARS-CoV-2 was determined on either the QuantStudio 5 (Applied Biosystems) or Lightcycler LC480II (Roche) with the 2019-nCoV CDC assay (<https://www.fda.gov/media/134922/download>).

Viral transport medium from positive swabs (stored at 4°C) was collected for samples run on the Roche Cobas and Hologic Panther Fusion systems and in all other cases excess RNA was collected (stored at 4°C and collected within 4 days for samples tested by the AusDiagnostics assay, while all other RNA samples were initially frozen and thawed for collection). Excess positive SARS-CoV-2 inactivated swab samples (200 µl viral transport medium from nose and throat swabs inactivated in 200 µl Zymo DNA/RNA shield and 800 µl Zymo viral DNA/RNA buffer) were collected from the Cytology and the Microbiology Department and SARS-CoV-2-positive RNA extracts (~20 µl) were collected from the Microbiology Department and the BCRE as part of the COG-UK Consortium project (PHE Research Ethics and Governance Group R and D ref. no NR0195), with full details in Tables S1–S3. For inactivated swab samples, RNA was extracted using the Quick DNA/RNA Viral Magbead kit from step 2 of the DNA/RNA purification protocol (Zymo - <https://tinyurl.com/y2lqoneq>).

SARS-CoV-2-positive samples were transferred to the Quadram Institute Bioscience for sequencing. The lower cycle threshold (Ct) or take-off value produced by the SARS-CoV-2 assays in the Roche, AusDiagnostics, Altona Diagnostics and CDC assays were used to determine whether samples needed to be diluted for sequencing according to the ARTIC protocol [for AusDiagnostics results, 13 was added to the take-off value to generate an approximate Ct value – this is because 15 cycles of PCR are performed before a dilution step and a further 35 cycles of nested PCR (the take-off value is determined in the nested PCR)]. The SARS-CoV-2 assay in the Hologic Panther does not provide a take-off or Ct value but rather a combined fluorescence signal for both targets in relative fluorescence units (RLUs), and therefore all samples tested by the Hologic Panther were processed undiluted in the ARTIC protocol.

### Sequencing using ARTIC SARS-CoV-2 multiplex tiling PCR

cDNA and multiplex PCRs were prepared following the ARTIC nCoV-2019 sequencing protocol v2 [16]. Dilutions of RNA were prepared when necessary based on Ct values following ARTIC protocol guidelines. V3 CoV-2 primers ([https://github.com/artic-network/artic-ncov2019/tree/master/primer\\_schemes/nCoV-2019/V3](https://github.com/artic-network/artic-ncov2019/tree/master/primer_schemes/nCoV-2019/V3)) were used to



perform the multiplex PCR for SARS-CoV-2 according to the ARTIC protocol [16] with minor changes. Due to variable Ct values, all RNA samples used in the two ARTIC multiplex PCRs were run for 35 cycles. Odd and even PCRs were pooled and cleaned using a 1× SPRI bead clean with KAPA Pure Beads (Roche Catalogue No. 07983298001), according to the manufacturer's instructions. PCR products were eluted in 30 µl of 10 mM Tris-HCl buffer, pH 7.5, and cDNA was quantified using the QuantiFluor ONE dsDNA System (Promega). Libraries were prepared for sequencing on the Illumina or Nanopore platform and sequenced as described previously [17].

## Sequence analysis

Raw reads were demultiplexed using bcl2fastq (v2.20) (Illumina) allowing for zero mismatches in the dual barcodes to produce FASTQ files. The reads were used to generate a consensus sequence for each sample using an open source pipeline adapted from <https://github.com/connor-lab/ncov2019-artic-nf> (<https://github.com/quadram-institute-bioscience/ncov2019-artic-nf/tree/qib>). Briefly, read adapters were trimmed using TrimGalore (<https://github.com/FelixKrueger/TrimGalore>) and aligned to the Wuhan Hu-1 reference genome (accession MN908947.3) using BWA-MEM (v0.7.17) [18]; ARTIC amplicons were masked and a consensus was built using iVar (v1.2) with primary parameters 'ivar consensus -m 10 -q 20 t 0.75' [19].

## Quality control

Samples were prepared and sequenced in 96-well plates with one cDNA-negative control per plate and one RNA extraction-negative control, where applicable. Contaminated samples were removed from analysis (Table S2). The COG-UK consortium defines a consensus sequence as passing COG-UK basic quality control (QC) if: >50% of the genome is covered by confident calls or there is at least one contiguous sequence of more than 10 000 bases; and no evidence of contamination in the negative control. A confident call is defined as having 10× depth of coverage. If the coverage falls below these thresholds, the bases are masked with the character N indicating the base at that position is unknown or not available. Low-quality variants are also masked with Ns. The QC threshold for inclusion in the public database GISAID (Global Initiative on Sharing All Influenza Data) is higher, requiring that >90% of the genome is covered by confident calls and that there is no evidence of contamination. The COG-UK quality control criteria were used as the minimum requirements for lineage and phylogenetic analysis. Although we did not use a homoplasy-based test for batch effects as described previously [20, 21], a visual inspection of the consensus sequences for common SNPs across lineage boundaries for each batch is undertaken using snipit (<https://github.com/aineniambh/snipit>). Additionally civet reports (<https://github.com/COG-UK/civet>) are generated for each batch, with genomes from other labs included for context and are manually inspected. These reports identify for each

generated sequence the closest sample in the COG-UK and GISAID databases, where in many cases identical samples sequenced by other labs are present. This helps to monitor for lineage altering artefacts.

## Data availability

All consensus sequences were deposited in GISAID [22] if they met its minimum QC threshold. All raw sequence data and metadata [23] were deposited in the European Nucleotide Archive (ENA) [24]. In both cases this happened soon after sequence generation, facilitated through COG-UK, and using MRC CLIMB [25].

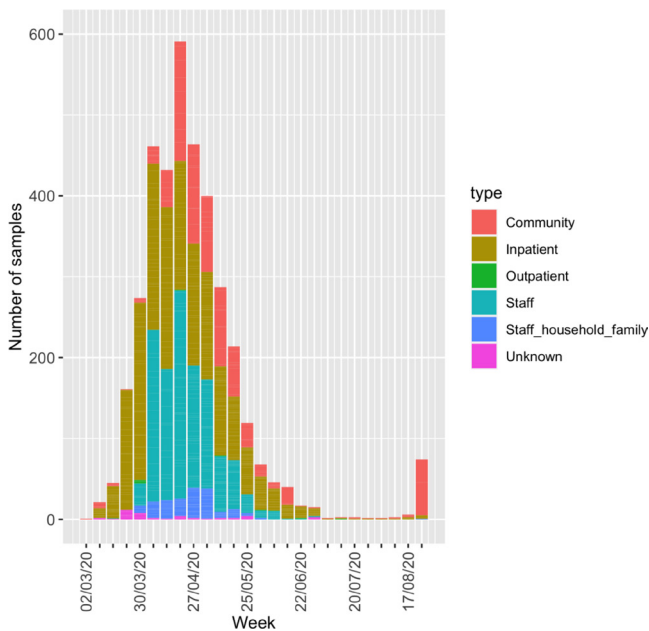
## Clustering and phylogenetic analysis

Lineages [13] assigned to each consensus genome were determined using Pangolin (<https://github.com/cov-lineages/pangolin>), which is run routinely by the Rambaut group over SARS-CoV-2 consensus sequences deposited on MRC CLIMB [25]. Global lineages are identifiers given to actively spreading lineages, are defined using a phylogenetic framework (<https://github.com/COG-UK/grapevine>) and often represent distinct introductions into new territories or regions, taking the form B.1.2.3 (see Rambaut et al. [13] for full details). UK lineages represent the subsequent spread within the UK, taking the form UK1234 and providing an identifier for a cluster for a given phylogeny. Unlike global lineages, however, UK lineage identifiers are not consistent between phylogenies. In this paper, a *sublineage* is defined as a set of samples within a lineage which share a common history within this lineage. For example, if all samples within a lineage are represented on a phylogenetic tree, a sublineage would describe all samples derived from a single internal node. Only samples that passed COG-UK QC were considered for lineage assignment (>50% of genome reconstructed), and only samples with more than 90% of non-ambiguous aligned sites (i.e. not N) were included in the upstream phylogenetic pipeline.

This upstream analysis at MRC CLIMB (<https://github.com/COG-UK/grapevine>) generates both the phylogenetic tree, aligned sequences and lineages (global and local) for most data. When adding more samples to the phylogenetic tree (as in the case studies below), we use IQ-Tree2 [26], first constraining the tree search to the upstream tree to create an initial tree, and then through unconstrained optimization using default parameters and the HKY model [27] with gamma heterogeneity [28] in both cases. Sequences are aligned with MAFFT using Wuhan Hu-1 (accession MN908947.3) as a reference genome. Phylogenetic trees were visualized using ggtree [29].

## Epidemiological analyses

Epidemiological analyses of outbreaks presented in the results were instigated and overseen by clinicians within the NHS or by public health bodies. Genome sequencing, bioinformatics analysis and genomic epidemiology were performed by Quadram Institute Bioscience with limited



**Fig. 1.** Total number of positive samples in the region per week, broken down by type. Not all of these were available for sequencing and a single individual may have been sampled multiple times. Staff (key workers) include healthcare workers and essential workers, such as police officers.

anonymized metadata as part of the COG-UK consortium; patient-identifiable data were retained by the hospitals or public health bodies.

### Longitudinal sampling and reinfection

For each case, the Norwich Research Park Biorepository (part of NNUH) anonymously linked all instances where cases were sequenced longitudinally (140) and provided the information to Quadram Institute Bioscience for analysis. The UK lineages were extracted for each case with multiple samples using precomputed lineages from COG-UK. Consensus genomes which did not yield high enough quality genomes to compute a lineage were excluded. Cases which had more than two high-quality samples were validated to ensure the lineages were the same. Where differences were identified, all consensus genomes for the case were extracted into a FASTA file, and the differences compared to the Wuhan Hu-1 reference were noted using SNP-sites [30] (version 2.3.3).

### Care facility outbreak analysis

An initial list of SARS-CoV-2 samples associated with a single care facility was provided by NNUH to Quadram Institute Bioscience. The UK lineages were identified for each sample using precomputed lineages from COG-UK. All other samples with the same UK lineage in the COG-UK dataset were identified and a phylogenetic tree was computed using IQ-Tree2 under an HKY+G model, as described above [26] (version 2.0.6). All samples in a sublineage associated with Norfolk were identified. The mutations defining this sublineage were

calculated using SNP-sites [30] (version 2.3.3), with the Wuhan Hu-1 reference as the base. A phylogenetic tree of the sublineage was calculated before first removing singleton mutations, most of which were C→T/U SNPs, and are markers of RNA degradation. Anonymized care facility sample metadata were added to the sublineage, with the data visualized in Phandango [31] and the relatedness of the samples and care facilities was visually confirmed. Hospital admission and discharge data for the residents was analysed solely by Public Health England co-authors and an anonymized summary was provided for this paper to maintain patient confidentiality.

### Hospital outbreak analysis

A list of SARS-CoV-2 samples associated with a hospital were provided by Ipswich Hospital to Quadram Institute Bioscience. The lineages were identified for each sample using precomputed lineages from COG-UK. The frequency of each lineage was identified for Ipswich Hospital. SARS-CoV-2 samples collected by hospitals (NNUH, QEH and JPUH) and care organizations in the same region as Ipswich Hospital had the frequencies of their lineages calculated in a similar fashion, providing context.

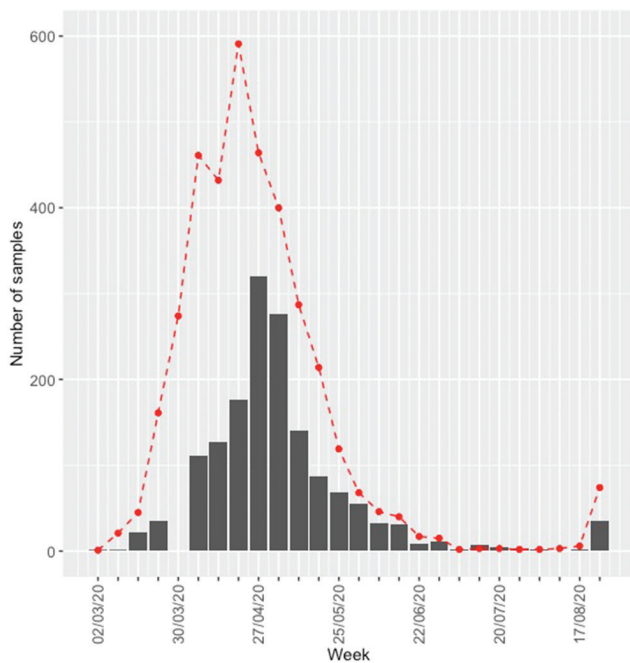
### Food processing facility outbreak analysis

SARS-CoV-2-positive samples were sent to Quadram Institute Bioscience. Samples were prepared, sequenced using an Oxford Nanopore Technologies MinION and bioinformatically analysed, all within 24 h of sample receipt. Consensus genomes were provided to Civet (<https://github.com/COG-UK/civet>), which assigned lineages to each genome. SNPs defining the sublineage of the outbreak were manually identified from the Civet report. The global lineage where the samples fell into was analysed further. All public samples from the global lineage which were publicly accessible through GISAID were identified. The countries of origin and collection dates were noted. A phylogenetic tree of all samples in this lineage from May onwards was created using IQ-Tree2 as described above. The SNPs defining the sublineages were calculated using SNP-sites (version 2.3.3).

## RESULTS

### Samples

The first reported case in the Norfolk region was on 6 March 2020 from a returning traveller; by 31 August 2020, there were 3225 cases identified by NNUH from Norfolk and surrounding areas from a total of 3751 SARS-CoV-2-positive clinical samples (some cases were sampled multiple times). Of these, 1565 samples (41.7%) were sequenced and analysed, from 1376 cases (42.6%). This represents approximately 172 SARS-CoV-2-positive samples sequenced per 100,000 population. The sequenced cases were broken down by locality and age group (Table S4). For cases sampled multiple times, the earliest collection date of a SARS-CoV-2-positive sample was used for sequence analysis. These samples were collected in the East of England, predominantly from cases providing an address in Norfolk (Fig. S1). The samples came from



**Fig. 2.** Number of positive samples sequenced at Quadram Institute Bioscience (grey bars) over time compared with the number of samples collected (red dashed line). Project sample collection only officially began on 8 April 2020 and not all samples taken before this time were available for sequencing. However, archived samples for March were sequenced and are represented in this figure.

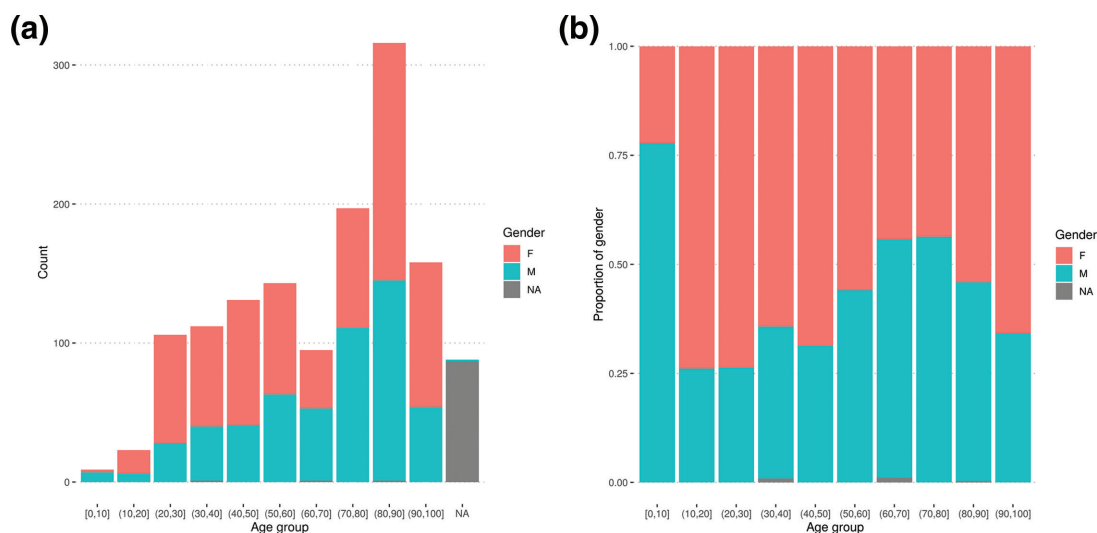
individuals in the community (20.7%,  $n=285$ ), inpatients (40.6%,  $n=559$ ) and outpatients (0.3%,  $n=4$ ) at hospitals, and staff (key workers) (23.8%,  $n=328$ ) and their families (4.7%,  $n=65$ ) (Fig. 1). Inpatients represented a mixture of patients newly admitted to the hospital (with or without COVID-19

symptoms), and existing patients with possible nosocomial SARS-CoV-2 infections. As testing was extended to more groups, so the regions from which samples were collected also changed (Fig. S2).

The number of positive samples in the Norfolk region peaked at the end of April 2020; specifically, the number of SARS-CoV-2-positive samples peaked in the week of 20–26 April 2020, with 591 positive samples (Fig. 2). The peak month was April with 1992 positive samples, followed by May with 1188 positive samples. These numbers include a small number of repeat samples. In July only 10 new positive cases were reported, before rising to 79 in August. More than 60 of the August cases were related to a food processing facility outbreak.

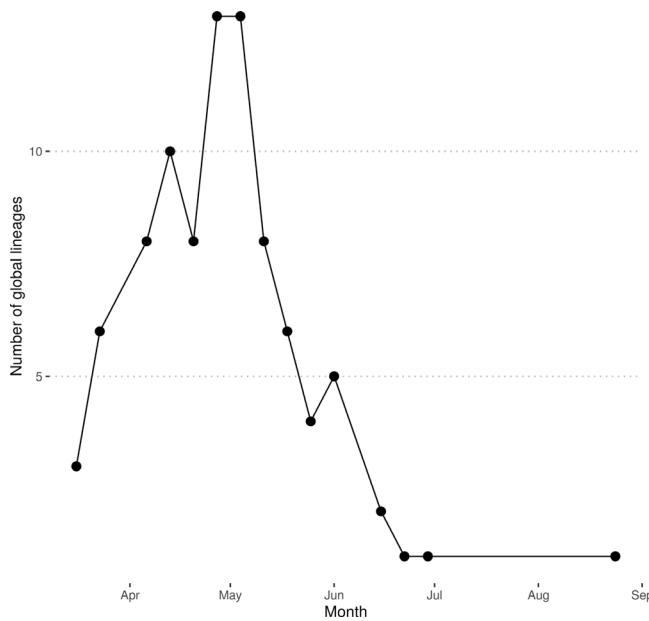
Proportionally, the number of SARS-CoV-2-positive samples that were sequenced followed the same trend as the total number of positive samples, peaking in the week of 27 April to 3 May, with genomes from 320 cases being sequenced (Fig. 2). Although project sample collection for sequencing officially began on 8 April 2020, and no samples were sequenced from the period 27 March to 7 April 2020, 59 archived samples from March were available and were sequenced. Overall the number of genomes sequenced does reflect the number of positive cases in the region and we can confidently conclude that the peak period was April/May 2020, in this region.

The number of positive cases sequenced was greatest in older individuals, with the largest number of samples ( $n=316$ , 36.5%) being from cases aged 80–90 years (Fig. 3a). Just nine samples originated from cases under 10 years of age. Females were also significantly over-represented in the dataset (Fig. 3b), accounting for 57% ( $n=741$  out of 1286,  $P<0.001$ , one-proportion z-test) of cases. Virtually all samples made available for sequencing were sequenced, with no selection criteria based on patient data or quality cut-offs; while we



**Fig. 3.** Age and sex of positive cases in absolute values (a) and scaled as proportions of all positive cases sequenced (b). Data from one person over 100 years of age are not included. NA, not available.





**Fig. 4.** The weekly number of co-occurring global lineages, excluding lineages that were only found in a single sample.

cannot rule out differential submission rates of testing facilities representing variations in patient populations, we thus generally expect the sequenced positive cases to reflect patient characteristics in the underlying population of total positive cases.

## Genomic analysis

### Genome quality control

Samples received for sequencing varied substantially in viral load. For our samples the Ct correlated well with the percentage of bases missing in the reconstructed consensus genome (Fig. S3) with a substantial reduction in genome completeness for samples with a Ct above 32. Virtually all diagnostic positives were sequenced, irrespective of Ct value, to avoid underrepresenting patients with low viral loads. In a small number of instances [32] samples were not available for sequencing as there was insufficient material or the samples could not be found.

We identified potential biases in the quality of the genomes sequenced from samples. The sex and age of the cases (for those genomes which had the relevant associated metadata) were evaluated against four QC categories: not sequenced, failed all QC, passed basic QC and passed high-quality QC (Fig. S4). There were significantly more genomes from females in three QC categories (62, 63 and 55% female in genomes failing all QC, passing basic QC and passing high-quality QC, respectively, one-proportion z-test,  $P < 0.001$ ), but not significantly different in the not sequenced category (66% female,  $P = 0.052$ ). With respect to an individual's age, the mean age of individuals contributing samples that failed all QC was significantly higher (mean age 70.0 years vs. 64.9, 65.7 and 58.8 years of age of individuals contributing genomes passing

basic QC, passing high-quality QC and not sequenced, respectively,  $P < 0.05$ , pairwise Kruskal–Wallis rank sum tests adjusted for multiple hypothesis testing).

Completeness of consensus genomes is related to the Ct of the input samples (Figs S5 and S6) and three temperature-sensitive ARTIC PCR primer dropout areas were visible [33] visible. A further three performed poorly at higher Ct values, due to reduced amplification efficiency, variation in which is to be expected in a large amplicon pool. Overall, the ARTIC protocol was robust and sensitive. Above Ct 32 the completeness of the genomes recovered did begin to tail off and there was a substantial, largely random, drop off above Ct 35, i.e. there was no consistency in the primer pairs that performed well or poorly in the multiplex when there were  $< 10$  genome copies present.

In total, 901 samples (65.4%) passed the GISAID QC criteria ( $\geq 90\%$  genome completeness), 120 samples (8.7%) passed the COG-UK QC criteria only ( $\geq 50\%$  genome completeness and 357 samples (25.9%) failed ( $< 50\%$  genome completeness).

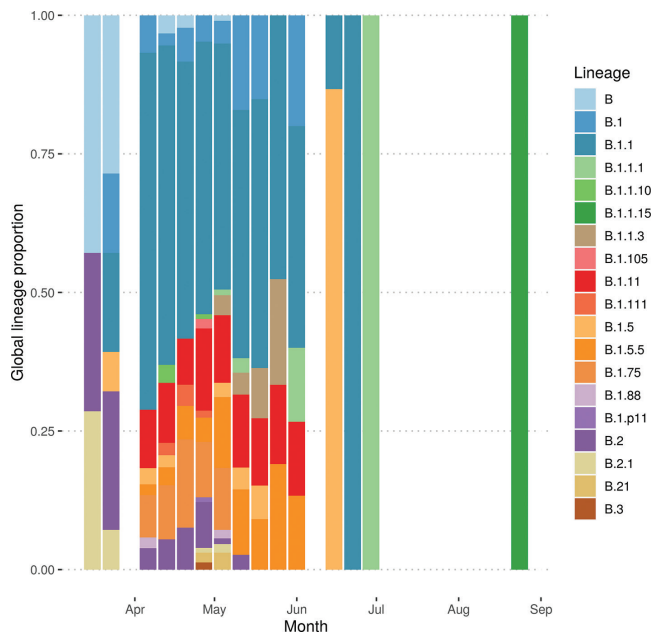
### Global lineages

As samples were collected over a 6-month period, the median number of SNPs per genome increased every month compared with the Wuhan Hu-1 reference (accession MN908947.3). When only considering high-quality consensus genomes (Fig. S7) it increased from six SNPs in March to 16 SNPs in August. The evolutionary rate was estimated to be  $\sim 2$  SNPs per month.

The maximum number of co-occurring global lineages in a given week was 13 for the period 27 April to 10 May 2020, approximately 5–6 weeks after the UK government instituted a lockdown (23 March 2020) (Fig. 4). This rapidly reduced as the number of samples dropped.

When considering the number and proportion of co-occurring global lineages every week it is apparent that, during the peak (April/May), some global lineages became extinct and were replaced by new lineages, which rapidly increased in abundance (Fig. 5). The transient nature of lineages is exemplified in lineages B.1.1.3 and B.2. Lineage B.1.1.3 was first identified in the week beginning 4 May 2020 (seven cases); it persisted for a month before becoming extinct within the region. Lineage B.2 was observed much earlier, in the week beginning 16 March 2020 (two cases), peaked 1 month later (19 cases) and was last observed on 11 May 2020 (two cases). The small number of lineages and low case numbers in July and August coincided with a low instance of disease in the region, including a negligible level of hospital admissions for COVID-19.

The global lineage nomenclature system of Rambaut *et al.* [13] provides a flexible and consistent naming scheme for genomically detectable introductions of SARS-CoV-2 into new locations; there have been 1035 global lineages assigned in this scheme. A total of 26 of these global lineages were observed in our data, 20 of which were observed in more than one sample. All of the global lineages present in more than one sample were from lineage B, which is



**Fig. 5.** The proportion of samples represented by each lineage per week, excluding lineages represented by a single sample. Sample collection for sequencing began on 8 April 2020, and no archival samples were available for sequencing for the period 26 March to 7 April 2020. Samples from July and August were primarily repeated sampling of the same case or with only a single individual in a lineage.

predominantly associated with European and UK outbreaks ([https://github.com/hCoV-2019/lineages/blob/master/lineages/data/lineage\\_descriptions\\_2020-05-19.txt](https://github.com/hCoV-2019/lineages/blob/master/lineages/data/lineage_descriptions_2020-05-19.txt)). A number of lineages were over- or under-represented in Norfolk samples compared with the UK and globally (Table 1). For example: B.1.1.15 was represented in 3.38% ( $n=35$ ) of Norfolk genomes compared with 0.97% ( $n=60$ ) globally (only 25 from outside Norfolk); B.1.11 was represented in 10.72% of Norfolk genomes compared with 1.34% globally; B.1.5.5 was represented in 6.47% of Norfolk genomes compared with 0.12% globally; and B.1.75 was represented in 8.12% of Norfolk genomes compared with 0.25% globally.

The most commonly sequenced global lineage was B.1.1, which was present in 464 samples (44.8%) in this dataset and 19424 globally (31.4%); this lineage has three defining SNPs at 28881 G→A, 28882 G→A and 28883 G→C and is closely associated with the European expansion of SARS-CoV-2.

Norfolk samples were set in context as part of the COG-UK phylogenetic pipeline (7 September 2020) using a phylogenetic tree based on all publicly released genome sequences (Fig. 6). Building phylogenetic trees with incomplete genomes is challenging [20, 34] in light of the low diversity of SARS-CoV-2 during the study period, with the genomes in the 50–90% completeness range having a higher potential for misalignment and phylogenetic misplacement. Overall genomes from the Norfolk region represented a random sampling of co-occurring global lineages within the UK as

a whole. Some major global lineages are under-represented in Norfolk such as B.1 and others are over-represented, such as B.1.1, when compared with the UK samples (Table 1).

### UK lineages

Global lineages were further subdivided into UK lineages, to identify ongoing transmission and evolution within the UK. The numbers assigned to UK lineages are subject to change and must be recalculated for all genomes with each phylogenetic reanalysis. Thus, the numbers reported here are for a single phylogenetic analysis. Stable cluster identification and nomenclature for SARS-CoV-2 is currently an open problem, and thus all analyses in this paper were relative to a single snapshot, with consistent algorithm/database/software versions, phylogenetic analysis and sample sets.

There were 100 UK lineages detected in the dataset, 49 of which were present in two or more cases. The number of co-occurring UK lineages peaked at 20 in the week of 27 April 2020, approximately 5 weeks after the UK nationwide lockdown began; thereafter the number dropped to a single lineage in July and August (Fig. S8). The proportions of samples with particular lineages varied week to week (Fig. S9), with the most common UK lineage being UK5, which was present in 324 cases; this is also the most commonly identified lineage in the UK (<https://microreact.org/project/cogconsortium-2020-09-02/f5aa0bdd/>). The next most common UK lineage was UK2913, which was present in 113 samples; this was a sublineage associated with care facilities in the region around Norwich city (detailed later).

### Examination of the D614G mutation in the spike protein

There is evidence that a mutation in the spike protein of SARS-CoV-2 (an amino acid change from D to G at position 614; D614G) increases infectivity of a pseudotype virus *in vitro* in cells; this is associated with an observed increase in viral loads in patients [35]. Overall, in the Norfolk dataset, 89.4% ( $n=819$ ) of samples had the D614G mutation while only 10.6% ( $n=97$ ) had the wild type (Fig. 7). The relative proportion of the two genotypes differed over time. In March, 66.6% ( $n=24$ ) of samples contained the wild type and 33.3% ( $n=12$ ) contained the D614G mutation. In April the proportion of genomes that were wild type had reduced to 10.7% ( $n=47$ ) while those with the D614G mutation were dominant at 89.3% ( $n=392$ ). In May the proportion of genomes that were wild type had reduced to 5.5% ( $n=22$ ) compared with 94.4% ( $n=374$ ) of genomes having the D614G mutation (Fig. 7).

### Epidemiological case studies

#### Longitudinal sampling and reinfection

In August 2020 Quadram Institute Bioscience, at the request of the Microbiology Department at the NNUH, evaluated the genomes present in multiple longitudinal samples taken from the same case over extended periods of time during infection. The aim was to determine whether they were infected by the same lineage or different lineages, the latter



**Table 1.** Summary of the prevalence of global lineages represented by more than one sample in the Norfolk dataset compared with the public set of genomes from GISAID and COG-UK (accessed 28 August 2020)

Lineage	Norfolk		UK <sup>a</sup>		Global <sup>b</sup>	
	Samples	Per cent	Samples	Per cent	Samples	Per cent
B	21	2.03	1021	2.59	1831	2.97
B.1	63	6.09	5918	14.99*	13354	21.63*
B.1.1	464	44.83	15654	39.65*	19424	31.46*
B.1.1.1	11	1.06	3658	9.26*	3933	6.37*
B.1.1.10	6	0.58	295	0.75	324	0.52
B.1.1.15	35	3.38	45	0.11*	60	0.97*
B.1.1.3	19	1.84	100	0.25*	100	0.16*
B.1.1.7	2	0.19	157	0.40	173	0.28
B.1.105	7	0.68	90	0.23	90	0.15*
B.1.11	111	10.72	793	2.01*	825	1.34*
B.1.111	11	1.06	11	0.03*	18	0.03*
B.1.5	33	3.19	997	2.53	1996	3.23
B.1.5.5	67	6.47	73	0.18*	75	0.12*
B.1.75	84	8.12	152	0.38*	156	0.25*
B.1.88	6	0.58	39	0.1*	43	0.07*
B.2	51	4.93	1343	3.40	1792	2.90*
B.2.1	13	1.26	2004	5.08*	2153	3.49*
B.21	12	1.16	26	0.07*	27	0.04*
B.3	9	0.87	1000	2.53*	1118	1.81

a, In this column those values marked with an asterisk are significantly more/less prevalent in Norfolk samples than UK samples ( $P < 0.05$  Bonferroni corrected for multiple hypothesis testing).

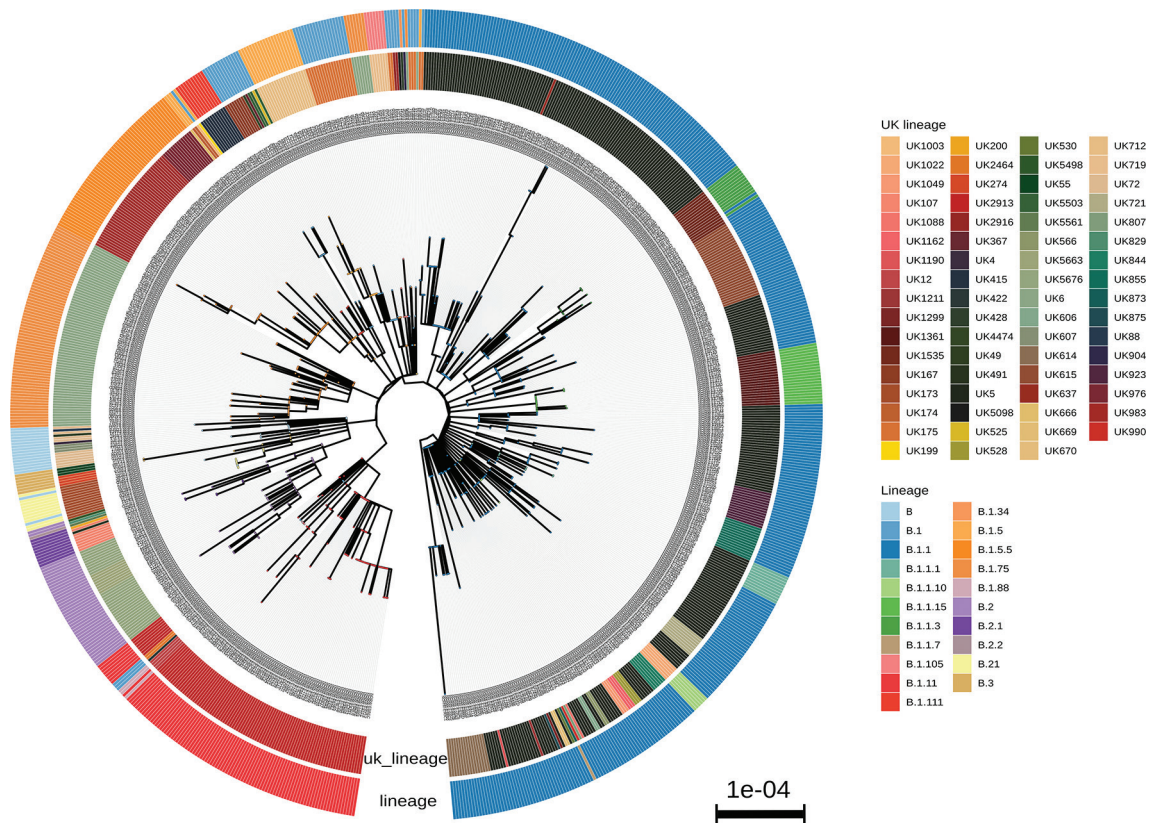
†‡b, In this column those values marked with an asterisk are significantly more/less prevalent in Norfolk samples than in global samples ( $P < 0.05$  Bonferroni corrected for multiple hypothesis testing).

indicating potential reinfection. Longitudinal SARS-CoV-2-positive samples were identified from 140 cases (Table S3), each with between two and six samples, 88.5% of whom were hospital inpatients at some point during their illness. The median time span of the sampling was 13 days, with a mean of 16.2 days. The longest time span was 71 days, with 22 cases having a time span greater than 28 days. The clinical outcomes were not available for analysis. Only samples with different collection dates from the same individual were considered. We limited cases to those with high-quality consensus genomes (passing GISAID QC) in two or more of the sample time points; this resulted in a series of longitudinal genome samples from 42 cases; each series had two to four samples. In every example the lineage remained the same between samples from the same individual, with the exception of the linked samples NORW-ED449 and NORW-ECD30; this was because NORW-ECD30 had nine IUPAC [36] symbols for ‘partially’ ambiguous bases, which are likely to be due to differences in viral load between the

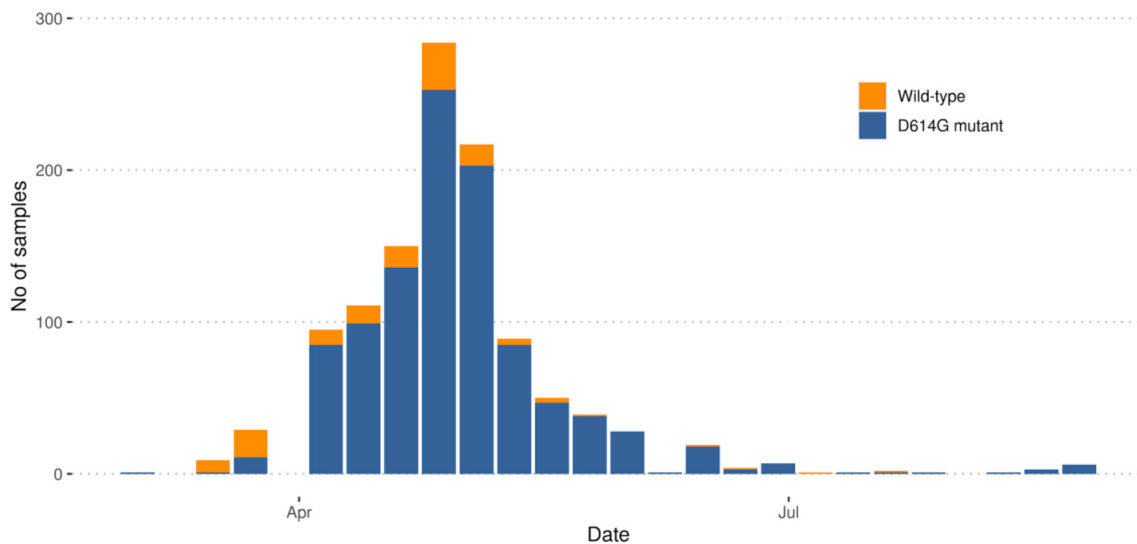
original samples (Ct 17 vs. Ct 27). These results suggest that there is no evidence of reinfection in any of the individuals for which a series of positive samples had been taken.

### Hospital outbreak analysis

In August 2020 Quadram Institute Bioscience and Ipswich Hospital used the genome data from a set of 31 samples from hospital patients to determine whether these samples represented a single nosocomial outbreak or whether they were unrelated and the result of community transmission. The 31 positive samples were collected between 6 March and 28 August 2020; 80.6% ( $n=25$ ) were from patients over the age of 65 years. From these, 18 yielded genome sequences of sufficient quality to assign lineages. A total of six global lineages and eight UK lineages were observed in the samples, with the most commonly observed ( $n=5$ ) being UK5, which is also the most commonly observed lineage within the UK. This number of co-occurring lineages



**Fig. 6.** Phylogenetic tree of the Norfolk genomes sequenced in this study. The phylogenetic tree was estimated as part of the COG-UK phylogenetic pipeline (7 September 2020). The inner circle represents the UK lineages assigned to each sequence, while the outer circle shows their equivalent global lineages. Only high-quality samples are included (837 sequences). The tree is scaled by nucleotide substitutions per site.



**Fig. 7.** Weekly numbers of genomes with the wild type genome (D614) or the D614G mutant in the spike protein.

indicated that there was not a single large nosocomial outbreak at this location.

### Care facility outbreak analysis

In June 2020 the Microbiology Department at NNUH and Quadram Institute Bioscience evaluated an outbreak at a care facility in the Norwich region using SARS-CoV-2 sequenced genomes from the dataset in this paper. The analysis undertaken indicated probable intra-care facility transmission, corresponding to a discrete sublineage circulating in the care sector as opposed to the wider community. It revealed that 14 out of 15 genomes from cases had the same UK lineage, UK2913, over a sustained period of time (8 April to 1 June 2020). An analysis of this lineage in all COG-UK data ( $n=395$ ) revealed that it represented a distinct sublineage in the Norwich region of Norfolk, defined by a single synonymous mutation (A→G) at position 24232 in the S gene; this mutation was not found in any other COG-UK lineages. Most of the cases with this sublineage were >80 years of age and concentrated in distinct areas in the Norwich region. Analysis also confirmed that the samples were predominantly collected from six care facilities. There were 89 cases sequenced in the Norwich region with this sublineage and 76 of these were known to be patients ( $n=64$ ) of care facilities or healthcare workers in those facilities ( $n=9$ ) and their families ( $n=3$ ). Links could not be established for 13 cases who tested positive for this sublineage. This sublineage had not been observed previously in community testing and the last new positive patient with this sublineage was on 1 June 2020. As it has not been seen in 3 months this sublineage is now regarded as extinct.

An analysis was undertaken to understand the role of hospital discharges in this sublineage. Of the residents in this sublineage, 12 had a hospital admission, of whom two were admitted twice to three hospital trusts. Six had a community-acquired infection, testing positive within 7 days of admission, three were inconclusive due to missing data, one had a probable hospital-acquired infection and tested positive within 7 days of discharge, and two had a definite hospital-acquired infection (<https://www.gov.uk/government/publications/wuhan-novel-coronavirus-infection-prevention-and-control/epidemiological-definitions-of-outbreaks-and-clusters-in-particular-settings>). All residents with a hospital-acquired infection had a test prior to discharge, suggesting the package of infection prevention and control (IPC) measures were being followed. In the time period covered by the study, patients required a test prior to discharge to a care facility, and a positive test did not preclude them from returning to the care facility, but rather that enhanced IPC and isolation measures needed to be taken for a designated period of time. Given that some of this cohort of patients tested positive in May with community-acquired infections, a number of weeks after the adult social care IPC Department of Health and Social Care measures were announced, suggests that these measures may not have been sufficient.

On examination of all genome sequences obtained from a town with two of these care facilities, we found 70 samples

were positive for SARS-CoV-2, with 52 of those yielding genome sequences of sufficient quality to assign a lineage. Thirty-seven samples (71%) were associated with the care facilities and were the UK2913 lineage. The remaining 15 samples in the town came from 13 different lineages, indicating that the number of co-occurring lineages within the care facilities did not reflect the number of co-occurring lineages within the wider locality.

### Food processing facility outbreak analysis

In August 2020, 35 positive samples from workers at a food processing facility were rapidly sequenced (<24h) as part of an outbreak analysis by Quadram Institute Bioscience, the Specialist Virology Centre at NNUH, and Norfolk County Council Test and Trace and Public Health England. A total of 27 yielded genomes of sufficient quality to assign lineages. All genomes fell within the same cluster with a global lineage of B.1.1.15 and a UK lineage of UK1361, confirming that these genomes were related to each other. The consensus sequences of the genomes in this outbreak had three SNPs defining a sublineage: U/T→C at position 7783, G→U/T at position 25552 and G→U/T at position 28221. This sublineage had not been observed before in the Norfolk data, COG-UK data or global GISAID data.

Looking more widely at the lineage, it was only observed in samples from the UK, with seven other samples identified by community sequencing in July (15–27 July 2020). However, these seven samples lacked the three sublineage-defining SNPs found in the outbreak samples, and had two additional defining SNPs of their own (G→U/T at position 12067, G→U/T at position 28086). This indicates that these two sets of samples had a common ancestor and, based on an evolutionary rate of ~2 SNPs per month, this ancestor might have occurred 1–2 months earlier. Ancestors of lineage B.1.1.15 were only observed in Portugal in May and June (GISAID, accessed 27 August 2020).

## DISCUSSION

Here, we have used intensive whole genome sequencing of SARS-CoV-2 samples in a single geographical area to investigate the evolution and transmission of the virus in this region. The average age of the Norfolk population is significantly higher than that of England as a whole; 24.5% of residents are aged 65 years or older compared with 18.4% for England as a whole (<https://www.norfolkinsight.org.uk/population/>, accessed 31 May 2020). The largest hospital in the region is the Norfolk and Norwich University Hospitals NHS Foundation Trust (1200 beds), serving a population of around one million patients from Norfolk and neighbouring counties, supported by a network of smaller hospitals.

There was a lower incidence of SARS-CoV-2 in the Norfolk region compared with England as a whole, with the proportion of cases testing positive for SARS-CoV-2 in Norfolk at 363.2 per 100000 compared with 573.9 per 100000 in England as a whole (accessed 14 September 2020, <https://coronavirus.data.gov.uk/>). This was also reflected in the number of deaths



that occurred due to infection within 28 days of a positive diagnosis; for Norfolk this was 43.8 per 100000 population compared with 65.7 per 100000 for England as a whole, which is significantly below average despite having an older, more vulnerable population.

The sequencing data represented a rate of 172 sequenced genomes per 100000 in the Norfolk population, which corresponded to 113.8 cases for which high-quality genome sequences were available for evaluation per 100000 in the population. Specifically, we evaluated high-quality genome sequence for 31.3% of all cases that tested positive for SARS-CoV-2 in Norfolk. However, as the samples primarily came from healthcare settings, they captured cases with the most severe symptoms rather than asymptomatic community cases. Only positive tests from the population with clinical need (primarily hospital and care facilities), and key workers and their families were available for sequencing. Community testing was done at large regional 'Lighthouse labs' capable of processing hundreds of thousands of samples per day. Most community testing required the case to have COVID-19 symptoms, with the exception of population-level surveillance, which made up a small proportion of positive cases. Of the sequenced community samples during the study period just 25 were attributed to Norfolk and lacked metadata, and thus were not included in analysis.

The highest number of positive samples in this region occurred at the end of April/beginning of May 2020, approximately 6 weeks after the UK instituted a nationwide lockdown (23 March 2020). Thereafter, the number of positive cases we sequenced dropped substantially as the impact of the lockdown and social distancing began to reduce transmission. With the exception of the food processing facility outbreak, by August 2020, only two new cases with positive samples were detected.

Analysis of the demographic metadata associated with positive samples indicated cases were more likely to be older and female. This skewed distribution in relation to age and sex is likely to be due to the directed use of diagnostic testing to symptomatic cases during the peak of the pandemic; this approach was driven by global shortages in reagents and testing capacity. Thus, vulnerable elderly cases were more likely to be tested during the peak of the pandemic, and were more likely to be female as they have a longer life expectancy. Thereafter, when testing was opened up to key workers, they were predominantly female healthcare workers who make up 77% of the NHS workforce [37].

Viral RNA loads in individuals, as measured by PCR, were correlated strongly with the percentage of the genome that could be reconstructed from the sequencing data. This could be due to individual variation in host factors, disease stage [38] or quality of the sample material [32]. Phylogenetically useful genomes where the Ct was below 32 (more than ~100 viral copies) were routinely sequenced but there was a substantial tail off from Ct 35. These results must be interpreted cautiously in terms of transmission potential, as it is not known whether an individual is infectious at low viral

RNA loads; it is possible that the positive results from high Ct samples are due to detection of residual RNA from a past infection. However, this does demonstrate that sequencing can produce usable information from samples containing the wide range of viral RNA loads likely to be encountered during SARS-CoV-2 infections.

There is evidence that a mutation in the spike protein with an amino acid change of D to G at position 614 (D614G) increases the transmissibility of the virus, which is associated with an increased viral load in mutant-infected cases [15, 35]. This has been observed in the UK and globally [15], potentially indicating that a more transmissible strain is now in circulation. As seen in Fig. 7, in March 2020, during the early part of the pandemic in the Norfolk region, the wild type was found in the majority of samples (66.6%), although the overall number of sequenced samples was low. However there was a rapid fall in the proportion of samples with the wild type in April (10.7%), and by May 95% of samples had the D614G mutation. This corresponds to a changing landscape of global lineages in the region, with an increasing dominance of B.1 lineages and sublineages, and a decreasing occurrence of other B (non-B.1) lineages.

The information provided by these sequences allowed an examination of the overall genetic variation within SARS-CoV-2 circulating in Norfolk and comparison with other regions. The number of co-occurring global lineages was similar to the range found within the UK as a whole [7], Europe [6] and beyond [39]. The notable exception was the lack of lineage A samples within the region, with only two being observed. This indicates that most of the lineages that entered the region did not come directly from China; rather, they are estimated to have predominantly come from Europe or within the UK. In the region, 23.2% ( $n=26$ ) of all global lineages were observed out of 112 lineages that have been defined to date [13]. This variation shows that genomically distinct lineages have expanded worldwide, with different distributions taking hold in different settings (see Microreact [40] <https://microreact.org/project/cogconsortium-2020-09-02/f5aa0bdd/>). The B.1.11 samples in our dataset were specifically associated with care facilities in Norfolk.

These data demonstrate a substantial number of co-occurring global lineages within one small region, indicating multiple concurrent introductions and their subsequent spread. This places a lower bound on the number of independent introductions to the region at 26, but it is likely to be substantially higher as not all COVID-19 infections were identified, tested and sequenced. As case numbers rose during the course of the pandemic, more lineages were identified, with a peak in the number of co-occurring lineages around 5–6 weeks after the UK instituted a national lockdown. Thereafter, the number of lineages dropped substantially, with many rapidly becoming extinct in the region, providing further evidence that lockdown measures break transmission.

By subdividing the global lineages into standardized UK lineages, a finer resolution of viral genomic relatedness was obtained and allows us to make more detailed comparisons.



As we observed with global lineages, the UK lineages provided further evidence for substantial viral genomic variability circulating in the region. There were 100 UK lineages observed in the region out of a total of 1725 lineages reported for the UK as a whole (5.7%). The dominant lineage in the UK (UK5) was also the dominant lineage in the region. Interestingly, the second most commonly observed lineage in the UK (UK1535), with 2269 (5.7%) samples and widespread circulation around the UK, was only observed 19 (1.8%) times in Norfolk (18 in this dataset, one sequenced by the University of Cambridge and none in Pillar 2 community testing for Norfolk).

One of the most important applications of these data was in epidemiological investigations to identify outbreaks. This was particularly important during the peak in April and May, when the number of new infections was high, providing the resolution required to distinguish between transmission clusters that would not have been otherwise possible. Our genomic data were informative in the following cases:

- (1) In a hospital setting, lineage information was used to differentiate nosocomial from community transmission. We sequenced 31 samples from Ipswich Hospital and found eight UK lineages; the most commonly observed (UK5) was also the most commonly observed in the UK. This demonstrated that a single large nosocomial outbreak had not occurred but that the patients in hospital had become infected in the community by circulating lineages.
- (2) Our data unexpectedly uncovered a sustained outbreak in six care facilities within the region. These data indicated probable intra-care facility transmission and is currently under further investigation. The outbreak was identified while looking at a common lineage within the region and noticing that the positive samples had mostly come from elderly people, suggesting a possible link to care facilities. Further investigation at NNUH identified six care facilities sharing a distinct sublineage primarily found only in these facilities. This sublineage was not detected in community testing (Pillar 2) at any point. Only two cases were a definite hospital acquired infection with one probable hospital-acquired infection. Examination of all genome sequences obtained from a town with two of these care facilities showed there were 13 different lineages circulating within the locality, but only a single lineage circulating in the care facilities.
- (3) Rapid genome sequencing was applied to an outbreak in a food processing facility in the region, yielding lineages for 27 genomes in 24 h. Sequencing this subset of samples (about 25% of positive cases detected) identified virtually identical genomes that were not found in the general community prior to the outbreak. This demonstrated within-factory transmission of a new lineage. An additional G→A mutation at position 20125 found in three of the outbreak genomes demonstrated that within the short period of time that the virus was circulating in the factory, some evolution with subsequent onward spread had occurred. As the lineage had not been seen

in the region previously, we are currently monitoring for transmission from the factory into the community. The factory sub-lineage is sufficiently novel to link it to the factory if identified anywhere in the country. Looking more widely at all publicly available data in GISAID, the only other country in which this global lineage (B.1.1.15) was recently observed was Portugal between 6 May and 14 June 2020 (the most recently observed genome from Portugal was from 22 June 2020). Genomics can be used to identify the possible source of a new introduction into a new area.

Samples could be broken down by the areas from which the samples were collected, representing towns or cities and their surrounds (Figs S2 and S10). Whilst UK5 was the commonest UK lineage in most urban areas in the Norfolk region, as it was nationally in the UK, there were other UK lineages that showed sustained persistence and spread within discrete communities from small geographical areas. UK2913 was primarily observed in Norwich and to the south and east of the city (South Norfolk, Waveney, Broadland, Great Yarmouth), nearly exclusively from clients of care facilities and healthcare workers at those facilities. Another UK lineage, UK721, was observed in seven community care residents (aged 78–92 years) and two healthcare workers in south-west Norfolk. One lineage, UK173, was observed only in one suburb of Norwich; it dominated for 1 month (13 April to 19 March 2020), then went extinct and has not been observed in Norfolk since. This indicates that lineages introduced into small urban areas may expand but do not necessarily spread more widely. The UK6 lineage was primarily observed in the Kings Lynn area, accounting for 90% (73/81) of UK6 samples in Norfolk. In contrast Norwich, 70 km away, recorded only a single case of UK6. These patterns are repeated throughout the dataset.

The number of co-occurring global and UK lineages in circulation amongst inpatients ( $n=559$ , 22 global, 71 UK) was also reflected in key workers and their families ( $n=394$ , 19 global, 49 UK). Sixteen UK lineages were observed in key workers and their families but were not observed in patients or in community care, though it must be noted that 10 of these lineages were only observed once. Four UK lineages (UK244, UK606, UK1049, UK1162) were observed three times each, and only in key workers and their families. As these lineages were not observed in hospitals or care facilities, it is likely that these key workers became infected in the community. It also indicates that these key workers were unlikely to have passed the virus to patients, providing evidence that infection control measures were effective in these cases. All UK lineages observed in household members of key workers were also seen in the key workers. Definitive confirmation of transmission, and the direction of transmission, cannot be inferred from the genome sequences due to the low evolutionary rate of the virus, which we have observed in our dataset as approximately two changes per month. Over time most lineages became dormant or extinct, with some expanding rapidly such as UK448 and then disappearing just as quickly over a 3-week period. Tracking lineages over time is an open challenge, with current methods requiring

continual reanalysis, which adds substantial complexity and a certain degree of uncertainty.

Caution must be exercised when making comparisons across different geographical regions, as the UK accounts for 64% (39483 out of 61740 as at 1 September 2020) of all publicly sequenced SARS-CoV-2 genomes. Out of 103 countries that have made SARS-CoV-2 genomes publicly available through GISAID, only the UK, Australia, Spain, India and the USA have sequenced more genomes than have been sequenced in Norfolk alone for this paper. The density of sequences when normalized for population size is 13 times greater in Norfolk than in the USA.

## CONCLUSION

We provide an in-depth examination of the genomic epidemiology of SARS-CoV-2 within a single geographical region covering the whole of the first wave of the pandemic from March to August 2020. We sequenced the genomes from 172 SARS-CoV-2-positive samples per 100000 population (1035 cases), representing 42.6% of all positive samples collected through the Microbiology Department at NNUH. From this, we identified 100 distinct lineages in the region, corresponding to multiple parallel introductions of the virus ( $n \geq 26$ ).

Dense sequencing of the virus provided actionable information for pandemic management, including: identifying a sublineage associated with care facilities, ruling out a large nosocomial outbreak in a hospital, showing no evidence of reinfection in longitudinal samples, and confirming an outbreak at a food processing facility while allowing for spillover into the community to be monitored. These achievements were only possible through the collaborative efforts of scientists (data and molecular), clinicians, data managers and epidemiologists. The large-scale genome sequencing of SARS-CoV-2-positive samples has provided valuable additional data for public health epidemiology in the Norfolk region, and will continue to help identify and untangle hidden transmission chains as the pandemic evolves.

### Funding information

The authors gratefully acknowledge the support of the Biotechnology and Biological Sciences Research Council (BBSRC); this research was funded by the BBSRC Institute Strategic Programme Microbes in the Food Chain BB/R012504/1 and its constituent projects BBS/E/F/000PR10348, BBS/E/F/000PR10349, BBS/E/F/000PR10351 and BBS/E/F/000PR10352. D.J.B., N.F.A., T.L.V. and A.J.P. were supported by the Quadram Institute Bioscience BBSRC-funded Core Capability Grant (project number BB/CCG1860/1). E.M.A. was funded by the BBSRC Institute Strategic Programme Gut Microbes and Health BB/R012490/1 and its constituent project(s) BBS/E/F/000PR10353 and BBS/E/F/000PR10356. The sequencing costs were funded by the COVID-19 Genomics UK (COG-UK) Consortium which is supported by funding from the Medical Research Council (MRC) part of UK Research and Innovation (UKRI), the National Institute of Health Research (NIHR) and Genome Research Limited, operating as the Wellcome Sanger Institute. The author(s) gratefully acknowledge the UKRI Biotechnology and Biological Sciences Research Council's (BBSRC) support of The Norwich Research Park Biorepository. L.G. was supported by a DART MRC iCASE and

Roche Diagnostics. A.P.T. was funded by Sara Borrell Research Grant CD018/0123 from ISCIII and co-financed by the European Development Regional Fund (A Way to Achieve Europe programme) and A.P.T.'s QIB internship was additionally funded by 'Ayuda de la SEIMC'. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

### Acknowledgements

Thanks to the COG-UK Consortium Study Group for their contributions. Thanks to Dr Judith Pell for her help and insightful comments on the manuscript. We gratefully acknowledge the submitters to GISAID; full details are listed in Supplementary Material 2.

### Author contributions

All authors have read the manuscript and consented to its publication. The study was designed and conceived by J.O.G., A.J.P. Paper writing was by A.J.P., A.E.M. Metadata analysis was by A.J.P., T.L.V., A.E.M., E.M. Bioinformatics analysis and informatics were performed by T.L.V., L.O.M., N.F.A., A.J.P. Sequencing and library preparation was performed by D.J.B., S.R., G.L.K., A.A., A.P.T., A.B., A.J.T., N.M.T., R.G., J.O.G. Clinical diagnostics and extractions were managed by A.K., S.D., R.P., N.E., E.M., L.C., L.B., R.D.D.A., C.B., undertook discharge analysis of care facility residents. Genomic epidemiological analyses performed by A.J.P., with oversight from M.C., L.S., S.D., E.M., L.B. Samples and metadata were collected by L.G., W.P., R.D., A.B., A.V.G., E.M.A., A.K., T.S., A.J.T. and M.D. and biobanked by R.S., RNA was extracted by A.B., A.J.T. Risk assessments were by G.L.K., R.A.K., J.W. Project management and oversight was by G.L.K., J.O.G., A.J.P., L.M., M.W., A.E.M., J.W. Funding for the project was secured by J.O.G., A.J.P.

### Conflicts of interest

L.G. received a partial support for his PhD from Roche. The other authors declare that there are no conflicts of interest.

### Ethical statement

The COVID-19 Genomics UK Consortium has been given approval by Public Health England's Research Ethics and Governance Group (PHE R and D Ref: NR0195). The studies were conducted as part of surveillance for COVID-19 infections under the auspices of Section 251 of the NHS Act 2006 and/or Regulation 3 of The Health Service (Control of Patient Information) Regulations 2002. They therefore did not require individual patient consent or ethical approval.

### References

- Huang C, Wang Y, Li X, Ren L, Zhao J, *et al.* Clinical features of patients infected with 2019 novel coronavirus in Wuhan, China. *The Lancet* 2020;395:497–506.
- Dong E, Du H, Gardner L. An interactive web-based dashboard to track COVID-19 in real time. *Lancet Infect Dis* 2020;20:533–534.
- Jin JM, Bai P, He W, Wu F, Liu XF, *et al.* Gender differences in patients With COVID-19: focus on severity and mortality. *Front Public Health* 2020;8.
- Bryant J, Chewapreecha C, Bentley SD. Developing insights into the mechanisms of evolution of bacterial pathogens from whole-genome sequences. *Future Microbiol* 2012;7:1283–1296.
- Lu R, Zhao X, Li J, Niu P, Yang B, *et al.* Genomic characterisation and epidemiology of 2019 novel coronavirus: implications for virus origins and receptor binding. *The Lancet* 2020;395:565–574.
- Alm E, Broberg EK, Connor T, Hodcroft EB, Komissarov AB, *et al.* Geographical and temporal distribution of SARS-CoV-2 clades in the WHO European Region, January to June 2020. *Euro Surveill* 2020;25:2001410.
- Filipe ADS, Shepherd J, Williams T, Hughes J, Aranday-Cortes E, *et al.* Genomic epidemiology of SARS-CoV-2 spread in Scotland highlights the role of European travel in COVID-19 emergence. *medRxiv* 2020;9.
- Oude Munnink BB, Nieuwenhuijse DF, Stein M, Á O, Haverkate M, *et al.* Rapid SARS-CoV-2 whole-genome sequencing and analysis for informed public health decision-making in the Netherlands. *Nat Med* 2020;26:1405–1410.
- Shu Y, McCauley J. GISAID: Global initiative on sharing all influenza data – from vision to reality. *Euro Surveill* 2017;22:30494.

10. Smith GJD, Vijaykrishna D, Bahl J, Lycett SJ, Worobey M, et al. Origins and evolutionary genomics of the 2009 swine-origin H1N1 influenza A epidemic. *Nature* 2009;459:1122–1125.
11. Page AJ, Alikhan N-, F, Carleton HA, Seemann T, Keane JA, et al. Comparison of classical multi-locus sequence typing software for next-generation sequencing data. *Microb Genomics* 2017;3.
12. Dallman T, Ashton P, Schafer U, Jironkin A, Painset A, et al. SnapperDB: a database solution for routine sequencing analysis of bacterial isolates. *Bioinformatics* 2018;34:3028–3029.
13. Rambaut A, Holmes EC, O'Toole Á, Hill V, McCrone JT, et al. A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology. *Nat Microbiol* 2020;5:1403–1407.
14. COG-UK. An integrated national scale SARS-CoV-2 genomic surveillance network. *Lancet Microbe* 2020.
15. Volz E, Hill V, McCrone JT, Price A, Jorgensen D, et al. Evaluating the effects of SARS-CoV-2 Spike mutation D614G on transmissibility and pathogenicity. *Cell* 2020.
16. Quick J. nCoV-2019 sequencing protocol v2. 2020. <https://www.protocols.io/view/ncov-2019-sequencing-protocol-v2-bdp7i5rn>
17. Baker DJ, Aydin A, Le-Viet T, Kay GL, Rudder S, et al. CoronaHiT: high-throughput sequencing of SARS-CoV-2 genomes. *Genome Med* 2021;13:21.
18. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *ArXiv13033997 Q-Bio* 2013.
19. Grubaugh ND, Gangavarapu K, Quick J, Matteson NL, De Jesus JG, et al. An amplicon-based sequencing framework for accurately measuring intrahost virus diversity using PrimalSeq and iVar. *Genome Biol* 2019;20:8.
20. Turakhia Y, Maio ND, Thornlow B, Gozashti L, Lanfear R, et al. Stability of SARS-CoV-2 phylogenies. *PLOS Genet* 2020;16:e1009175.
21. De Maio N, Walker C, Borges R, Weilguny L, Slodkowitz G, et al. Issues with SARS-CoV-2 sequencing data. *Virological* 2020.
22. Elbe S, Buckland-Merrett G. Data, disease and diplomacy: GISAID's innovative contribution to global health. *Glob Chall* 2017;1:33–46.
23. Griffiths EJ, Timme RE, Page AJ, Alikhan N-, F, Fornika D, et al. The PHA4GE SARS-CoV-2 contextual data specification for open genomic epidemiology. 2020. <https://www.preprints.org/manuscript/202008.0220/v1>
24. Cochrane G, Karsch-Mizrachi I, Takagi T. Sequence Database Collaboration IN. The international nucleotide sequence database collaboration. *Nucleic Acids Res* 2016;44:D48–50.
25. Connor TR, Loman NJ, Thompson S, Smith A, Southgate J, et al. CLIMB (the Cloud Infrastructure for Microbial Bioinformatics): an online resource for the medical microbiology community. *Microb Genomics* 2016;2:e000086.
26. Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, et al. IQ-TREE 2: new models and efficient methods for phylogenetic inference in the genomic Era. *Mol Biol Evol* 2020.
27. Hasegawa M, Kishino H, Yano T. Dating of the human-ape splitting by a molecular clock of mitochondrial DNA. *J Mol Evol* 1985;22:160–174.
28. Yang Z. Maximum likelihood phylogenetic estimation from DNA sequences with variable rates over sites: Approximate methods. *J Mol Evol* 1994;39:306–314.
29. Yu G. Using ggtree to visualize data on tree-like structures. *Curr Protoc Bioinforma* 2020;69:e96.
30. Page AJ, Taylor B, Delaney AJ, Soares J, Seemann T. Snp-sites: Rapid efficient extraction of Snps from multi-FASTA alignments. *Microb Genomics* 2016;2:e000056.
31. Hadfield J, Croucher NJ, Goater RJ, Abudahab K, Aanensen DM, et al. Phandango: an interactive viewer for bacterial population genomics. *Bioinformatics* 2018;34:292–293.
32. Rogers AA, Baumann RE, Borillo GA, Kagan RM, Batterman HJ, et al. Evaluation of transport media and specimen transport conditions for the detection of sars-cov-2 by use of real-time reverse transcription-pcr. *J Clin Microbiol* 2020;58.
33. Farr B, Rajan D, Betteridge E, Shirley L, Quail M, et al. COVID-19 ARTIC v3 Illumina library construction and sequencing protocol. 2020. <https://www.protocols.io/view/covid-19-artic-v3-illumina-library-construction-an-bgg3jvyn>
34. Morel B, Barbera P, Czech L, Bettisworth B, Hübner L, et al. Phylogenetic Analysis of SARS-CoV-2 Data Is Difficult. *Mol Biol Evol* 2020.
35. Korber B, Fischer WM, Gnanakaran S, Yoon H, Theiler J, et al. Tracking changes in sars-cov-2 spike: Evidence that d614g increases infectivity of the COVID-19 virus. *Cell* 2020;182:812–827.
36. IUPAC-IUB Comm. on Biochem. Nomenclature (CBN). Abbreviations and symbols for nucleic acids, polynucleotides, and their constituents. *Biochemistry* 1970;9:4022–4027.
37. NHS. Equality and diversity NHS trusts and CCGs March 2019 [Internet]. 2019. <https://digital.nhs.uk/data-and-information/publications/statistical/nhs-workforce-statistics/nhs-workforce-statistics---march-2019-provisional-statistics>
38. Zou L, Ruan F, Huang M, Liang L, Huang H, et al. SARS-COV-2 viral load in upper respiratory specimens of infected patients. *N Engl J Med* 2020;382:1177–1179.
39. van Dorp L, Acman M, Richard D, Shaw LP, Ford CE, et al. Emergence of genomic diversity and recurrent mutations in SARS-CoV-2. *Infect Genet Evol* 2020;83:104351.
40. Argimón S, Abudahab K, Goater RJE, Fedosejev A, Bhai J, et al. Microreact: Visualizing and sharing data for genomic epidemiology and phylogeography. *Microb Genom* 2016;2:e000093.

### Five reasons to publish your next article with a Microbiology Society journal

1. The Microbiology Society is a not-for-profit organization.
2. We offer fast and rigorous peer review – average time to first decision is 4–6 weeks.
3. Our journals have a global readership with subscriptions held in research institutions around the world.
4. 80% of our authors rate our submission process as 'excellent' or 'very good'.
5. Your article will be published on an interactive journal platform with advanced metrics.

Find out more and submit your article at [microbiologyresearch.org](https://microbiologyresearch.org).